

35.G2761



PATENT APPLICATION

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of:)
YASUO OKUTANI, ET AL.) Examiner: UNASSIGNED
Application No.: 09/818,607) Group Art Unit: 2641
Filed: March 28, 2001)
For: SPEECH SIGNAL)
PROCESSING APPARATUS)
AND METHOD, AND)
STORAGE MEDIUM) July 19, 2001

#3
7/27/01
MB

RECEIVED
JUL 20 2001
Technology Center 2600

Commissioner for Patents
Washington, D.C. 20231

CLAIM TO PRIORITY

Sir:

Applicants hereby claim priority under the International Convention and all rights
to which they are entitled under 35 U.S.C. § 119 based upon the following Japanese

Priority Application No.:

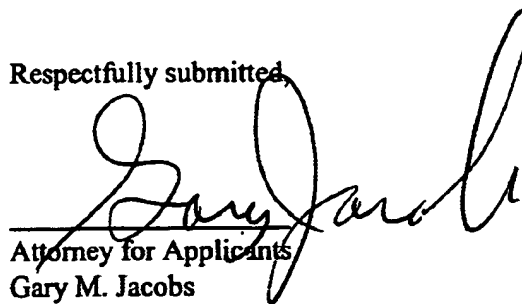
2000-099533 filed March 31, 2000.

A certified copy of the priority document is enclosed.

BEST AVAILABLE COPY

Applicants' undersigned attorney may be reached in our Washington, D.C. office by telephone at (202) 530-1010. All correspondence should continue to be directed to our below-listed address.

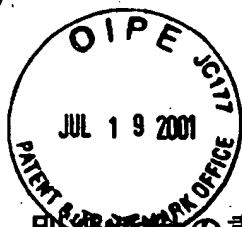
Respectfully submitted,

A large, stylized handwritten signature in black ink, appearing to read "Gary Jacobs".

Attorney for Applicants
Gary M. Jacobs
Registration No. 28,861

FITZPATRICK, CELLA, HARPER & SCINTO
30 Rockefeller Plaza
New York, New York 10112-3801
Facsimile No.: (212) 218-2200

DC_MAIN 59894 v 1



日 本 国 特 許 庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日

Date of Application:

2000年 3月31日

出 願 番 号

Application Number:

特願2000-099533

出 願 人

Applicant(s):

キヤノン株式会社

*Inventor: Yasuo Okutani, et al.
Applicant: 09/818, 607
Filed: 3/31/00*

RECEIVED

JUL 20 2001

Technology Center 2600

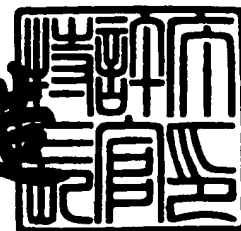
CERTIFIED COPY OF
PRIORITY DOCUMENT

BEST AVAILABLE COPY

2001年 4月20日

特 許 庁 長 官
Commissioner,
Japan Patent Office

及 川 耕 造



出証番号 出証特2001-3033160

【書類名】 特許願

【整理番号】 4172019

【提出日】 平成12年 3月31日

【あて先】 特許庁長官殿

【国際特許分類】 G01L 5/04

【発明の名称】 音声情報処理装置及びその方法と記憶媒体

【請求項の数】 19

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号 キヤノン株式会社
社内

【氏名】 奥谷 泰夫

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号 キヤノン株式会社
社内

【氏名】 小森 康弘

【特許出願人】

【識別番号】 000001007

【氏名又は名称】 キヤノン株式会社

【代理人】

【識別番号】 100076428

【弁理士】

【氏名又は名称】 大塚 康德

【電話番号】 03-5276-3241

【選任した代理人】

【識別番号】 100101306

【弁理士】

【氏名又は名称】 丸山 幸雄

【電話番号】 03-5276-3241

【選任した代理人】

【識別番号】 100115071

【弁理士】

【氏名又は名称】 大塚 康弘

【電話番号】 03-5276-3241

【手数料の表示】

【予納台帳番号】 003458

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 0001010

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 音声情報処理装置及びその方法と記憶媒体

【特許請求の範囲】

【請求項 1】 音声素片を所定の韻律情報に基づいて編集することによって生じる歪を求める歪出力手段と、

前記歪出力手段から出力された歪に基づいて、音声合成に使用する音声素片を選択する選択手段と、

前記選択手段により選択された音声素片に基づいて、音声合成を行う音声合成手段と、

を有することを特徴とする音声情報処理装置。

【請求項 2】 前記選択手段は、歪が最小となるように音声合成に使用する音声素片を選択することを特徴とする請求項 1 に記載の音声情報処理装置。

【請求項 3】 前記歪出力手段は、前記音声素片を他の音声素片と接続することによって生じる接続歪と前記音声素片を変形することによって生じる変形歪とに基づいて、前記歪を求めることを特徴とする請求項 1 又は 2 に記載の音声情報処理装置。

【請求項 4】 前記歪出力手段は、前記接続歪と、前記変形歪との重み付き和として前記歪を算出することを特徴とする請求項 3 に記載の音声情報処理装置。

【請求項 5】 前記歪出力手段は、ケプストラム距離を用いて前記接続歪を算出することを特徴とする請求項 3 又は 4 に記載の音声情報処理装置。

【請求項 6】 前記歪出力手段は、ケプストラム距離を用いて前記変形歪を算出することを特徴とする請求項 3 又は 4 に記載の音声情報処理装置。

【請求項 7】 前記歪出力手段は、前記変形歪を記憶したテーブルを有し、当該テーブルを参照して前記変形歪を決定することを特徴とする請求項 3 又は 4 に記載の音声情報処理装置。

【請求項 8】 前記歪出力手段は、前記接続歪を記憶したテーブルを有し、当該テーブルを参照して前記接続歪を決定することを特徴とする請求項 3 又は 4 に記載の音声情報処理装置。

【請求項 9】 テキストデータを入力する入力手段と、
前記テキストデータを言語解析する言語解析手段と、
前記言語解析手段により解析された結果に基づいて前記所定の韻律情報を生成する韻律情報生成手段と、
を更に備えることを特徴とする請求項 1 乃至 8 のいずれか 1 項に記載の音声情報処理装置。

【請求項 10】 音声素片を所定の韻律情報に基づいて編集することによって生じる歪を求める歪出力工程と、
前記歪出力工程で出力された歪に基づいて、音声合成に使用する音声素片を選択する選択工程と、
前記選択工程で選択された音声素片に基づいて、音声合成を行う音声合成工程と、
を有することを特徴とする音声情報処理方法。

【請求項 11】 前記選択工程では、前記歪が最小となるように音声合成に使用する音声素片を選択することを特徴とする請求項 10 に記載の音声情報処理方法。

【請求項 12】 前記歪出力工程では、前記音声素片を他の音声素片と接続することによって生じる接続歪と前記音声素片を変形することによって生じる変形歪とに基づいて、前記歪を求めることを特徴とする請求項 10 又は 11 に記載の音声情報処理方法。

【請求項 13】 前記歪出力工程では、前記接続歪と、前記変形歪との重み付き和として前記歪を算出することを特徴とする請求項 12 に記載の音声情報処理方法。

【請求項 14】 前記歪出力工程では、ケプストラム距離を用いて前記接続歪を算出することを特徴とする請求項 12 又は 13 に記載の音声情報処理方法。

【請求項 15】 前記歪出力工程では、ケプストラム距離を用いて前記変形歪を算出することを特徴とする請求項 12 又は 13 に記載の音声情報処理方法。

【請求項 16】 前記歪出力工程では、前記変形歪を記憶したテーブルを有し、当該テーブルを参照して前記変形歪を決定することを特徴とする請求項 12

又は 1 3 に記載の音声情報処理方法。

【請求項 1 7】 前記歪算出工程では、前記接続歪を記憶したテーブルを有し、当該テーブルを参照して前記接続歪を決定することを特徴とする請求項 1 2 又は 1 3 に記載の音声情報処理方法。

【請求項 1 8】 テキストデータを入力する入力工程と、
前記テキストデータを言語解析する言語解析工程と、
前記言語解析工程で解析された結果に基づいて前記所定の韻律情報を生成する韻律情報生成工程と、
を更に備えることを特徴とする請求項 1 0 乃至 1 7 のいずれか 1 項に記載の音声情報処理方法。

【請求項 1 9】 請求項 1 0 乃至 1 8 のいずれか 1 項に記載の方法を実行するプログラムを記憶したことを特徴とする、コンピュータにより読取り可能な記憶媒体。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は、音声素片を編集、接続して音声合成を行なう音声情報処理装置及びその方法と、その方法を実現するプログラムを記憶した記憶媒体に関するものである。

【0 0 0 2】

【従来の技術】

近年、テキストデータを入力し、そのテキストデータを言語解析してポーズ部分、無音時間の長さ、アクセントの生成などを行なって韻律情報を生成し、更に、その韻律情報に従って音声素片を記憶している素片辞書を検索し、対応する音声素片を読み出して音声合成する音声合成装置が知られている。

【0 0 0 3】

このような音声合成装置では、その読み出した音声素片を 1 ピッチ波形単位で複製、削除しながら、所望のピッチ間隔で貼り合わせて編集し（PSOLA：ピッチ同期波形重畳法）、それらの音声素片を接続する音声合成方式が主流となってい

る。

【 0 0 0 4 】

【発明が解決しようとする課題】

このような技術を利用して合成された音声には、音声素片を編集（変形）することによる歪（以下、変形歪）と、音声素片同士を接続することによって生じる歪（以下、接続歪）とが含まれ、これら 2 つの歪が合成音声の品質劣化を引き起こす大きな要因となっている。

【 0 0 0 5 】

本発明は上記従来例に鑑みてなされたもので、接続や変形に基づく歪の影響を小さくする音声情報処理装置及びその方法と、その方法を実現するプログラムを記憶した記憶媒体を提供することを目的とする。

【 0 0 0 6 】

【課題を解決するための手段】

上記目的を達成するために本発明の音声情報処理装置は以下のような構成を備える。即ち、

音声素片を所定の韻律情報に基づいて編集することによって生じる歪を求める歪出力手段と、

前記歪出力手段から出力された歪に基づいて、音声合成に使用する音声素片を選択する選択手段と、

前記選択手段により選択された音声素片に基づいて、音声合成を行う音声合成手段と、

を有することを特徴とする。

【 0 0 0 7 】

上記目的を達成するために本発明の音声情報処理方法は以下のような工程を備える。即ち、

音声素片を所定の韻律情報に基づいて編集することによって生じる歪を求める歪出力工程と、

前記歪出力工程で出力された歪に基づいて、音声合成に使用する音声素片を選択する選択工程と、

前記選択工程で選択された音声素片に基づいて、音声合成を行う音声合成工程と、
を有することを特徴とする。

【 0 0 0 8 】

【発明の実施の形態】

以下、添付図面を参照して本発明の好適な実施の形態を詳細に説明する。

【 0 0 0 9 】

〔実施の形態 1〕

図 1 は、本発明の実施の形態に係る音声合成装置のハードウェア構成を示すブロック図である。尚、本実施の形態では、一般的なパーソナルコンピュータを音声合成装置として用いる場合について説明するが、本発明は専用の音声合成装置であっても、また他の形態の装置であっても良い。

【 0 0 1 0 】

図 1 において、1 0 1 は制御メモリ（ROM）で、中央処理装置（CPU）1 0 2 で使用される各種制御データを記憶している。CPU 1 0 2 は、RAM 1 0 3 に記憶された制御プログラムを実行して、この装置全体の動作を制御している。1 0 3 はメモリ（RAM）で、CPU 1 0 2 による各種制御処理の実行時、ワークエリアとして使用されて各種データを一時的に保存するとともに、CPU 1 0 2 による各種処理の実行時、外部記憶装置 1 0 4 から制御プログラムをロードして記憶している。この外部記憶装置は、例えばハードディスク、CD-ROM 等を含んでいる。1 0 5 は D/A 変換器で、音声信号を示すデジタルデータが入力されると、これをアナログ信号に変換してスピーカ 1 0 9 に出力して音声を再生する。1 0 6 は入力部で、オペレータにより操作される、例えばキーボードや、マウス等のポインティングデバイスを備えている。1 0 7 は表示部で、例えば CRT や液晶等の表示器を有している。1 0 8 はバスで、これら各部を接続している。1 1 0 は音声合成ユニットである。

【 0 0 1 1 】

以上の構成において、本実施の形態の音声合成ユニット 1 1 0 を制御するための制御プログラムは外部記憶装置 1 0 4 からロードされて RAM 1 0 3 に記憶さ

れ、その制御プログラムで用いる各種データは制御メモリ101に記憶されている。これらのデータは、中央処理装置102の制御の下にバス108を通じて適宜メモリ103に取り込まれ、中央処理装置102による制御処理で使用される。D/A変換器105は、制御プログラムを実行することによって作成される音声波形データ（デジタル信号）をアナログ信号に変換してスピーカ109に出力する。

【0012】

図2は、本実施の形態に係る音声合成ユニット110の構成を示すブロック図である。

【0013】

図2において、201は入力部106や外部記憶装置104から任意のテキストデータを入力するテキスト入力部、202は解析辞書、203は言語解析部、204は韻律生成規則保持部、205は韻律生成部、206は素片辞書である音声素片保持部、207は音声素片選択部、208は音声素片編集・接続部、209は音声波形出力部である。

【0014】

以上の構成において、言語解析部203が、解析辞書202を参照して、テキスト入力部201から入力されるテキストの言語解析を行なう。こうして解析された結果が韻律生成部205に輸入される。韻律生成部205は、言語解析部203における解析結果と、韻律生成規則保持部204に保持されている韻律生成規則に関する情報とを基に音韻系列と韻律情報を生成して音声素片選択部207及び音声素片編集・接続部208に出力する。続いて、音声素片選択部207は、韻律生成部205から入力される韻律生成結果を用いて、音声素片保持部206に保持されている音声素片から対応する音声素片を選択する。音声素片編集・接続部208は、韻律生成部205から入力される韻律生成結果に従って、音声素片選択部207から出力される音声素片を編集及び接続して音声波形を生成する。こうして生成された音声波形は、音声波形出力部209で出力される。

【0015】

次に、上記構成を備えた本実施の形態の音声合成処理について説明する。

【0016】

図3は、本実施の形態に係る音声合成ユニット110における音声合成処理の流れを示すフローチャートである。

【0017】

まずステップS301で、テキスト入力部201は、文、文節、単語等の単位毎に、テキストデータを入力してステップS302に移る。ステップS302では、言語解析部203により当該テキストデータの言語解析を行う。次にステップS303に進み、音韻生成部205は、ステップS302で解析された結果と所定の韻律規則とに基づいて、音韻系列と韻律情報を生成する。次にステップS304に進み、各音韻毎に、ステップS303で得られた韻律情報と所定の音韻環境とに基づいて、音声素片選択部207が音声素片保持部206に登録されている音声素片を選択する。次にステップS305に進み、その選択された音声素片及びステップS303で生成された韻律情報とに基づいて、音声素片編集・接続部208により音声素片の編集および接続を行なってステップS306に進む。ステップS306では、音声素片編集・接続部208によって生成された音声波形を、音声波形出力部209が音声信号として出力する。このようにして、入力されたテキストに対応する音声出力されることになる。

【0018】

図4は、図3のステップS304（音声素片選択）の処理の詳細を示すフローチャートである。

【0019】

このステップS304では、音声素片同士の接続歪（後述する）と、音声素片の変形歪（後述する）とに基づいて決定される歪値（後述）に従って、動的計画法により、入力テキストデータの全体に互って歪値が最小となる音声素片系列を決定する。つまり、韻律生成部205が生成する音韻系列 P_n ($0 \leq n < N$) の先頭 ($n=0$) から順に処理することになる。まず最初の $n=0$ にセットし、ステップS401で、音韻系列の終端まで処理が終了していない場合、つまり $n < N$ の場合はステップS402に進み、 n 番目の音韻における音声素片の候補を音声素片保持部206から取り出し、それら音声素片の候補の個数を M_n としてス

テップ S403 に進む。ステップ S403 では、まず最初に $m=0$ にセットし、上記 n , m で特定される音声素片候補 $P_{n,m}$ ($0 \leq m < M_n$) に着目し、 n 番目の音韻における音声素片の候補の先頭 ($m=0$) から順に処理すして、候補の最後まで処理が終了していない場合、つまり $m < M_n$ の場合はステップ S404 に進むが、最後まで処理が終了した場合は、次の音韻の処理に移行するため $n = n + 1$ としてステップ S401 に戻る。ステップ S404 では、一つ前の音韻 P_{n-1} の各音声素片候補 $P_{n-1,k}$ ($0 \leq k < M_{n-1}$: M_{n-1} は一つ前の音韻 P_{n-1} の音声素片候補の数) と、候補 $P_{n,m}$ との間の歪値 $D_{k,m}$ をそれぞれ計算してステップ S405 に進む。ステップ S405 では、候補 $P_{n,m}$ に至るまでの歪値の総和の最小値である総和 $S_{n,m}$ を求める。この総和 $S_{n,m}$ は、次式で表現される。

【0020】

$$S_{n,m} = \min (S_{n-1,k} + D_{k,m}) \quad : 0 \leq k < M_{n-1}$$

この式において、 $\min()$ は、 k を“0”から“ M_{n-1} ”まで変化させた場合の最小値を意味する。また、そのときの k の値を $PRE_{n,m}$ として保持しておく。この $PRE_{n,m}$ は、候補 $P_{n,m}$ に至るまでの歪値の総和が最小となる経路を示し、ステップ S406 において最小歪経路を特定するために利用される。この候補 $P_{n,m}$ の総和 $S_{n,m}$ と $PRE_{n,m}$ が決定したら、次の音声素片候補に対する処理を行なうために $m = m + 1$ としてステップ S403 に戻る。

【0021】

こうしてステップ S401 で、最終である N 番目の音韻系列までの処理が終了するとステップ S406 に進み、総和 $S_{N-1,m}$ ($0 \leq m < M_n$) が最小となる候補 $P_{N-1,m}$ を特定し、そこから順次 $PRE_{n,m}$ を辿ることによって最小歪経路となる音声素片系列を特定する。こうして音声素片系列が特定されたら図3のステップ S305 に進んで、これら特定された音声素片の編集・接続を実行する。

【0022】

図5は、 n 番目の音韻（現在注目している音韻）の音声素片候補 $P_{n,1}$ における総和 $S_{n,1}$ の算出を模式的に示した図である。本実施の形態では、音韻の単位として *diphone* を採用した場合について記述する。

【0023】

図中、一つの円が音声素片の1候補 $P_{n,m}$ を示し、円内の数字が歪値の総和の最小値である総和 $S_{n,m}$ を示している。また矢印は、上述の $PRE_{n,m}$ を指す。また、四角で囲まれた数字は、音声素片候補 $P_{n,m}$ の歪値 $D_{k,m}$ を表わしている。

【0024】

次に、本実施の形態における歪値について説明する。

【0025】

ここでは、歪値 $D_{k,m}$ を接続歪 D_c と変形歪 D_t の重み付き和として定義する。

即ち、

$$D = w \times D_c + (1 - w) \times D_t \quad : (0 \leq w \leq 1)$$

ここで、重み係数 w は、予備実験など経験的に求める係数で、 $w = 0$ の場合は、歪値が変形歪 D_t のみで説明され、 $w = 1$ の場合は歪値が接続歪 D_c のみに依存することになる。

【0026】

図5では、音声素片候補 $P_{n,1}$ の一つ前の音韻の音声素片候補 $P_{n-1,250}$ との間の歪値 $D_{2,1}$ が“3”であり、音声素片候補 $P_{n-1,250}$ に至るまでの歪値の総和 $S_{n-1,2}$ が“8”であるため、経路51が $PRE_{n,1}$ として決定される。

【0027】

図6は、本実施の形態における接続歪 D_c の求め方を説明する図である。

【0028】

接続歪 D_c は、一つ前の音声素片と現在の音声素片との接続箇所において生じる歪で、本実施の形態では、ケプストラム距離を用いて表す。ここでは、音声素片境界が存在するフレーム60, 61（フレーム長5ミリ秒、分析窓幅25.6ミリ秒）と、それを挟む前後それぞれの2フレームからなる計5フレームを接続歪の算出対象とする。ケプストラムは、0次（パワー）～16次（パワー）までの計17次元とする。このケプストラムベクトルの各要素の差の絶対値の和を、現在注目している音声素片における接続歪とする。一つ前の音声素片における終端部のケプストラムベクトルの各要素を $C_{p i,j}$ （ i はフレーム番号で、 $i = 0$ が音声素片境界があるフレームである。 j はベクトルの要素番号を示す）、当該音声素片における始端部のケプストラムベクトルの各要素を $C_{c i,j}$ とすると、

現在注目している音声素片の接続歪 D_c は、

$$D_c = \sum \sum |C_{p\ i,j} - C_{c\ i,j}|$$

となる。ここで最初の \sum は $i = -2 \sim 2$ の総和を示し、次の \sum は $j = 0 \sim 16$ の総和を示している。

【0029】

図7は、本実施の形態に係る変形歪 D_c の求め方を説明する図である。

【0030】

ここでは、PSOLA法によりピッチ間隔を広げる場合について図示している。矢印はピッチマーク、点線は変形前と変形後のピッチ素片の対応関係を表わしている。本実施の形態では、各ピッチ素片（微細素片ともいう）の変形前後のケプストラム距離に基づいて変形歪を表す。具体的には、まず、変形後のあるピッチ素片（例えば70で示す）のピッチマーク71を中心にハニング窓72（窓長25.6ミリ秒）をかけ、そのピッチ素片70を周辺のピッチ素片を含めて切り出す。こうして切り出したピッチ素片70をケプストラム分析する。次に、ピッチマーク71に対応する変形前のピッチ素片73のピッチマーク74を中心にして同じ窓長のハニング窓75でピッチ素片を切り出し、変形後の場合と同様にしてケプストラムを求める。このようにして求めたケプストラム同士の距離を、着目しているピッチ素片70の変形歪として、変形後のピッチ素片とそれに対応する変形前のピッチ素片間の変形歪の総和をPSOLAで採用されるピッチ素片数 N_p で割った値を、その音声素片の変形歪とする。こうして求められる変形歪を式で記述すると以下のようなになる。

【0031】

$$D_t = \sum \sum |C_{org\ i,j} - C_{tar\ i,j}| / N_p$$

ここで最初の \sum は、 $i = 1$ から N までの総和を示し、次の \sum は $j = 0 \sim 16$ までの総和を示している。また $C_{tar\ i,j}$ は、変形後の i 番目のピッチ素片のケプストラムの j 次元目の要素を表わし、同様に、 $C_{org\ i,j}$ は、変形後に対応する変形前のピッチ素片のケプストラムの j 次元目の要素を表わしている。

【0032】

このように本実施の形態1によれば、各音声素片における接続歪及び変形歪を

求め、これら歪を基に重み付け計算を行なって各音声素片における歪値を求め、この歪値の総和が最小となる音声素片系列を特定して音声合成することにより、良好な音声合成結果を得ることができるという効果がある。

【0033】

〔実施の形態2〕

前述の実施の形態1では、音韻の単位としてdiphoneを用いる場合について記述したが、本発明はこれに限定されるものではなく、音素や半diphoneなどを単位としてもよい。半diphoneとは、diphoneを音素境界で2つに分割したもののことである。

【0034】

図8は、半diphoneを単位とした場合の概念図である。この半diphoneを単位とした場合のメリットについて簡単に説明する。任意のテキストを合成する場合、素片辞書は、全種類のdiphoneを用意しておく必要がある。これに対して、半diphoneを単位とした場合は、足りない半diphoneを別の半diphoneで代替できる。例えば、半diphoneの「/a.b.0/(diphone a.bの左側)」の代わりに「/a.n.0/」を利用しても、音質の劣化を少なくして良好に音声を再生できる。これにより、素片辞書のサイズをより小さくできる。

【0035】

〔実施の形態3〕

前述の実施の形態1及び2では、音韻の単位としてdiphoneや音素や半diphoneを用いる場合について説明したが、本発明はこれに限定されるものではなく、これらを混合して用いてもよい。例えば、利用頻度が高い音韻については、diphoneを単位として、利用頻度が低い音韻については、2つの半diphoneを用いて表現するようにしても良い。

【0036】

図9は、音声素片単位を混合した場合の一例を示した図で、ここでは音韻「o.w」がdiphoneで表され、その前後の音韻は半diphoneで表されている。

【0037】

〔実施の形態4〕

実施の形態3において、元のデータベース中で連続する場所から取り出されたかどうかの情報をもち、連続していた場合は、半diphoneの組を仮想的にdiphoneとして扱うようにしてもよい。つまり、元のデータベース中で連続するということは接続歪が“0”であるため、この場合には変形歪だけを考慮すればよいことになり計算量を大幅に軽減できる。

【0038】

図10は、この様子を表わした概念図である。図中の線上の数字は接続歪を表している。

【0039】

図10において、1100で示される半diphoneの組は、元のデータベース中で連続する場所から取り出されたものであり、その接続歪みは“0”に一義的に決定されている。また1101で示された半diphoneの組は、元のデータベース中で連続する場所から取り出されたものではないため、それぞれに対して接続歪みが計算される。

【0040】

[実施の形態5]

上述の実施の形態1では、動的計画法を、1単位のテキストデータから得られた音韻系列全体に対して適用する場合について説明したが、本発明はこれに限定されるものではない。例えば、ポーズや無音部分までを一つの区間として音韻系列を分割し、各区間毎に動的計画法を実行してもよい。尚、ここで言う無音部分とは、p,t,kなどの無音部分のことである。このようなポーズや無音部分では、接続歪が“0”であると考えられるため、このような分割が有効となる。これにより、区間ごとに適当な選択結果を得ることができるだけでなく、合成音声の生成に要する時間が短縮できる。

【0041】

[実施の形態6]

前述の実施の形態1では、接続歪の計算にケプストラムを用いる場合について説明したが、本発明はこれに限定されるものではない。例えば、接続点の前後に互る波形の差分の和を用いて接続歪を求めてもよい。またスペクトル距離など

を用いて接続歪を求めてもよい。この場合、接続点はピッチマークに同期させるのが、より好ましい。

【0042】

〔実施の形態7〕

前述の実施の形態1では、接続歪の計算において、窓長、シフト長、ケプストラムの次数、フレーム数などを具体的数字を使って説明したが、本発明はこれに限定されるものではない。任意の窓長、シフト長、次数、フレーム数を使って接続歪を算出してもよい。

【0043】

〔実施の形態8〕

前述の実施の形態1では、接続歪の計算にケプストラムの次数ごとに差分を取ったものの総和を用いる場合について説明したが、本発明はこれに限定されるものではない。例えば、各次数を統計的性質などを使って正規化（正規化係数 r_j ）してもよい。この場合の接続歪 D_c は、

$$D_c = \sum \sum (r_j \times |C_{pre\ i,j} - C_{cur\ i,j}|)$$

となる。ここで、最初の \sum は $i = -2 \sim 2$ の総和を、次の \sum は $j = 0 \sim 16$ までの総和を示す。

【0044】

〔実施の形態9〕

実施の形態1では、ケプストラムの次数ごとの差分の絶対値をベースに接続歪の算出を行なう場合について説明したが、本発明はこれに限定されるものではない。例えば、差分の絶対値の累乗（累数が偶数の場合は絶対値でなくてもよい）をベースに接続歪の算出を行なってもよい。ここで累数を N とすると、接続歪 D_c は、

$$D_c = \sum \sum |C_{pre\ i,j} - C_{cur\ i,j}|^N$$

となる。ここで“ N ”は N 乗を示す。ここで N の値を大きくすることは、大きな差分について敏感になることを意味しているので、その結果、接続歪が平均的に小さくなるように働くことになる。

【0045】

〔実施の形態 10〕

前述の実施の形態 1 では、変形歪としてケプストラムを用いる場合について説明したが、本発明はこれに限定されるものではない。例えば、変形前後の一定区間の波形の差分の和を用いて変形歪を求めてもよい。また、スペクトル距離などを用いて変形歪としてもよい。

【0046】

〔実施の形態 11〕

前述の実施の形態 1 では、変形歪を波形から得られる情報を基に算出する場合について説明したが、本発明はこれに限定されるものではない。例えば、PSOLA によるピッチ素片の削除および複製の回数などを変形歪を算出する要素としてもよい。

【0047】

〔実施の形態 12〕

前述の実施の形態 1 では、音声合成時に音声素片を読み出す毎に接続歪を計算する場合について説明したが、本発明はこれに限定されるものではない。例えば、接続歪を予め計算しておき、テーブルとして保持しておいてもよい。

【0048】

図 11 は、diphone 「/a.r/」 と diphone 「/r.i/」 との間の接続歪を記憶したテーブルの一例を示す図である。ここでは縦軸に「/a.r/」の音声素片、横軸に「/r.i/」の音声素片をとっている。例えば、「/a.r/」の「id3」の音声素片と「/r.i/」の「id2」の音声素片との接続歪は“3.6”で表されている。このように接続可能な diphone 間の接続歪を全てテーブル化して用意することにより、音声素片同士の合成時の接続歪の算出がテーブルの参照だけで済むため、その計算量を大幅に軽減でき、算出時間を大幅に短縮できる。

【0049】

〔実施の形態 13〕

前述の実施の形態 1 では、音声合成時に、音声素片編集する毎に変形歪を計算する場合について説明したが、本発明はこれに限定されるものではない。例えば、変形歪を予め計算しておき、テーブルとして保持しておいてもよい。

【0050】

図12は、あるdiphoneを基本周波数と音韻時間長について変化させた場合の変形歪をテーブルで表した図である。

【0051】

図中、 μ は、そのdiphoneの統計的な平均値を示し、 σ は標準偏差である。具体的な表の作成方法としては、次のような作成方法が考えられる。まず、基本周波数と音韻時間長に関して統計的に平均値と分散を求める。次に、それらを基に（ $5 \times 5 =$ ）25通りの基本周波数と音韻時間長をターゲットとしてPSOLA法をそれぞれ適用し、テーブルの変形歪を一つずつ求めていけばよい。合成時は、ターゲットの基本周波数と音韻時間長が決まれば、テーブルの近傍の値で内挿（もしくは外挿）することによって、変形歪を推定することが可能である。

【0052】

図13は、合成時に変形歪を推定するための具体例を示した図である。

【0053】

図中、黒丸がターゲットの基本周波数と音韻時間長であり、このとき、各格子点の変形歪がテーブルからA, B, C, Dと求まっていると仮定すると、変形歪Dtは、以下の式により求めることができる。

$$Dt = \{A \cdot (1 - y) + C \cdot y\} \times (1 - x) + \{B \cdot (1 - y) + D \cdot y\} \times x$$

【0054】

〔実施の形態14〕

前述の実施の形態13では、変形歪テーブルの格子点として、そのdiphoneの統計的な平均値と標準偏差を基に 5×5 のテーブルを作成したが、本発明はこれに限定されるものではなく、任意の格子点を持つテーブルとしてもよい。また、格子点を平均値などに依らず決定的に与えてもよいものとする。例えば、韻律推定で推定されうる範囲を等分割するなどもよいものとする。

【0055】

〔実施の形態15〕

前述の実施の形態1では、接続歪と変形歪の重み和で歪を定量化する場合について説明したが本発明はこれに限定されるものではなく、接続歪と変形歪それぞ

れに閾値を設定しておき、どちらか一方でもその閾値を越えた場合はその音声素片が選択されないようにして、十分大きな歪の値を与えるようにしてもよい。

【 0 0 5 6 】

上記実施の形態においては、各部を同一の計算機上で構成する場合について説明したが本発明はこれに限定されるものではなく、例えばネットワーク上に分散した計算機や処理装置などに分かれて各部を構成してもよい。

【 0 0 5 7 】

上記実施の形態においては、各部を同一の計算機上で構成する場合について説明したが、これに限定されるものではなく、ネットワーク上に分散した計算機や処理装置などに分かれて各部を構成してもよい。

【 0 0 5 8 】

上記実施の形態においては、プログラムを制御メモリ（ROM）に保持する場合について説明したが、これに限定されるものではなく、外部記憶など任意の記憶媒体を用いて実現してもよい。また、同様の動作をする回路で実現してもよい。

【 0 0 5 9 】

なお本発明は、複数の機器から構成されるシステムに適用しても、1つの機器からなる装置に適用してもよい。前述した実施の形態の機能を実現するソフトウェアのプログラムコードを記録した記録媒体を、システム或いは装置に供給し、そのシステム或いは装置のコンピュータ（またはCPUやMPU）が記録媒体に格納されたプログラムコードを読み出し実行することによっても達成される。

【 0 0 6 0 】

この場合、記録媒体から読み出されたプログラムコード自体が前述した実施の形態の機能を実現することになり、そのプログラムコードを記録した記録媒体は本発明を構成することになる。このようなプログラムコードを供給するための記録媒体としては、例えば、フロッピーディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

【 0 0 6 1 】

また、コンピュータが読み出したプログラムコードを実行することにより、前述した実施の形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているOSなどが実際の処理の一部または全部を行ない、その処理によって前述した実施の形態の機能が実現される場合も含まれる。

【0062】

更に、記録媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行ない、その処理によって前述した実施の形態の機能が実現される場合も含まれる。

【0063】

以上説明したように本実施の形態によれば、音声合成に際して、音声素片を選択する際に接続歪と変形歪を基準とするようにしたため、音質の劣化を最小限に抑えた音声素片系列を求めて音声合成できる。

【0064】

【発明の効果】

以上説明したように本発明によれば、接続や変形に基づく歪の影響を小さくした音声合成することができる。

【図面の簡単な説明】

【図1】

本発明の実施の形態に係る音声合成装置のハードウェア構成を示すブロック図である。

【図2】

本発明の実施の形態1に係る音声合成装置の機能構成を示すブロック図である。

【図3】

本実施の形態に係る音声合成装置における処理の流れを示すフローチャートである。

【図 4】

図 3 のステップ S 3 0 4 の音素素片選択処理の詳細を示すフローチャートである。

【図 5】

n 番目の音韻の音声素片候補 $P_{n,1}$ における最小歪の総和 $S_{n,1}$ の算出を模式的に示した図である。

【図 6】

本発明の実施の形態に係る音声素片の接続歪を説明する図である。

【図 7】

本発明の実施の形態に係る音声素片の変形を説明する図である。

【図 8】

半diphoneを単位とした場合の概念図である。

【図 9】

本発明の実施の形態 3 に係る音声素片の単位をdiphoneと半diphoneとで混合した場合を説明する図である。

【図 1 0】

本発明の実施の形態 4 に係る音声素片の単位を取り出した半diphoneによって混合した例を示した図である。

【図 1 1】

本発明の実施の形態 1 2 に係るdiphoneの /a.r/ と /r.i/ 間の接続歪を決定するテーブル構成例を示す図である。

【図 1 2】

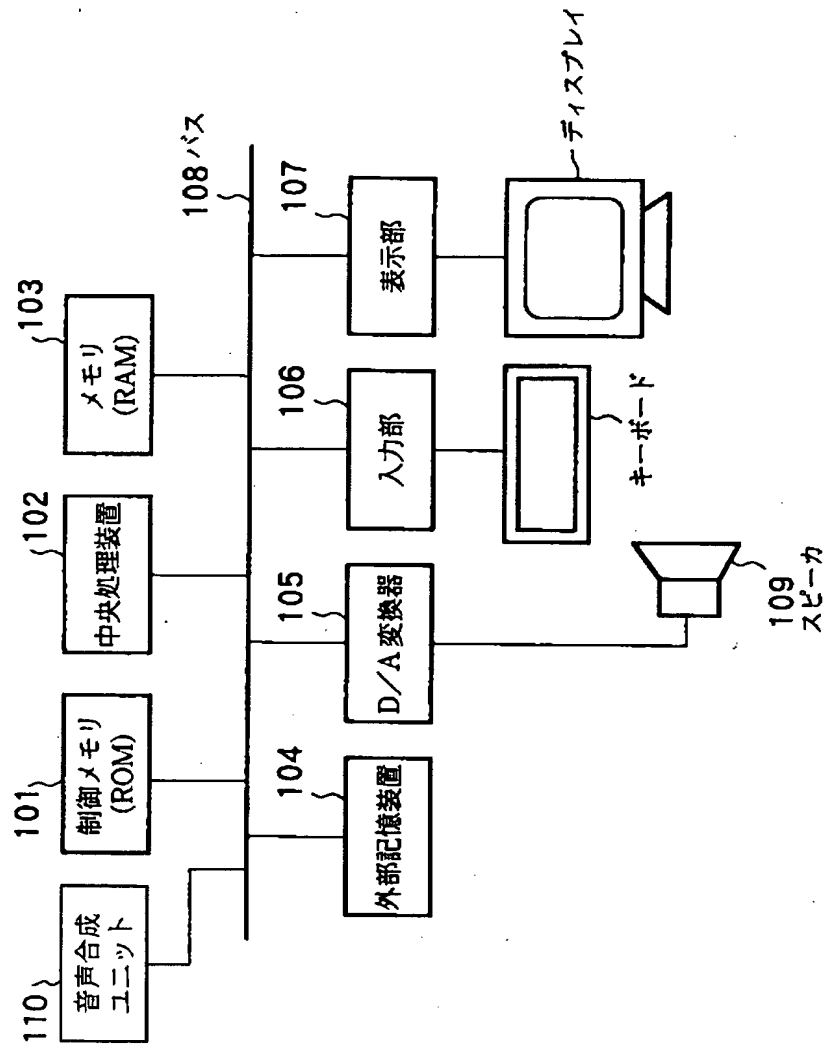
本発明の実施の形態 1 3 に係る変形歪を表わすテーブル例を示す図である。

【図 1 3】

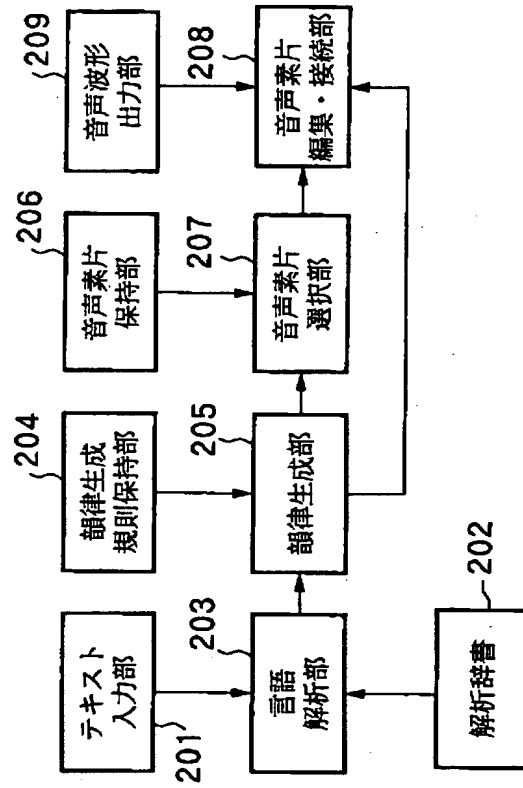
本発明の実施の形態 1 3 に係る変形歪を推定する具体例を示した図である。

【書類名】 図面

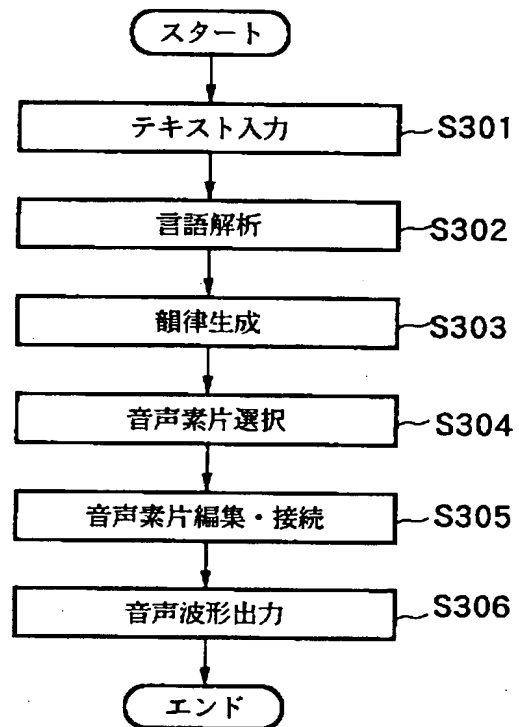
【図 1】



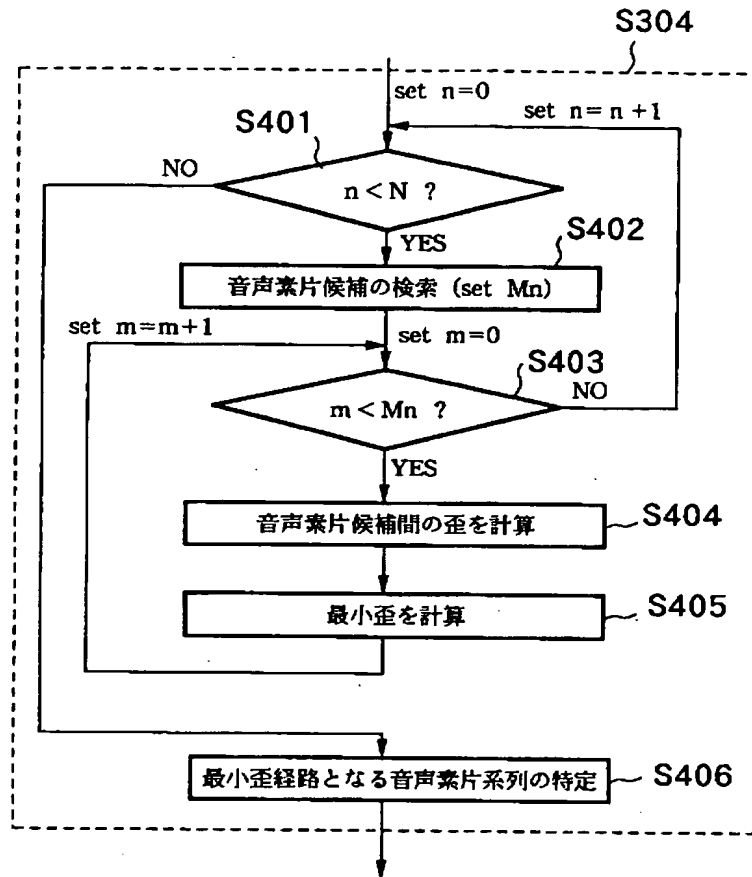
【図 2】



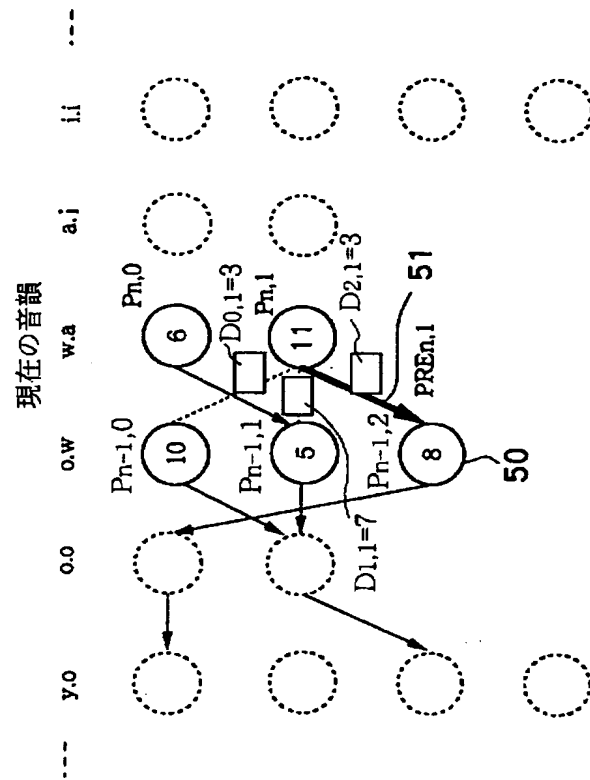
【図 3】



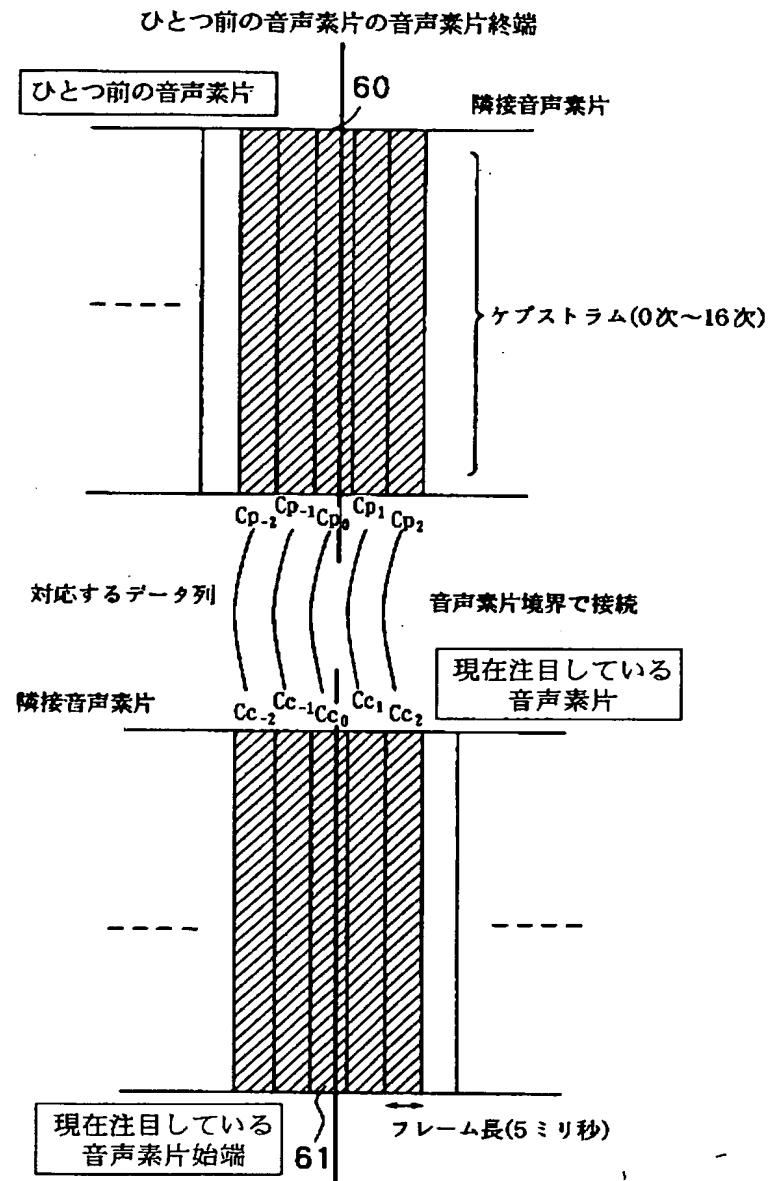
【図 4】



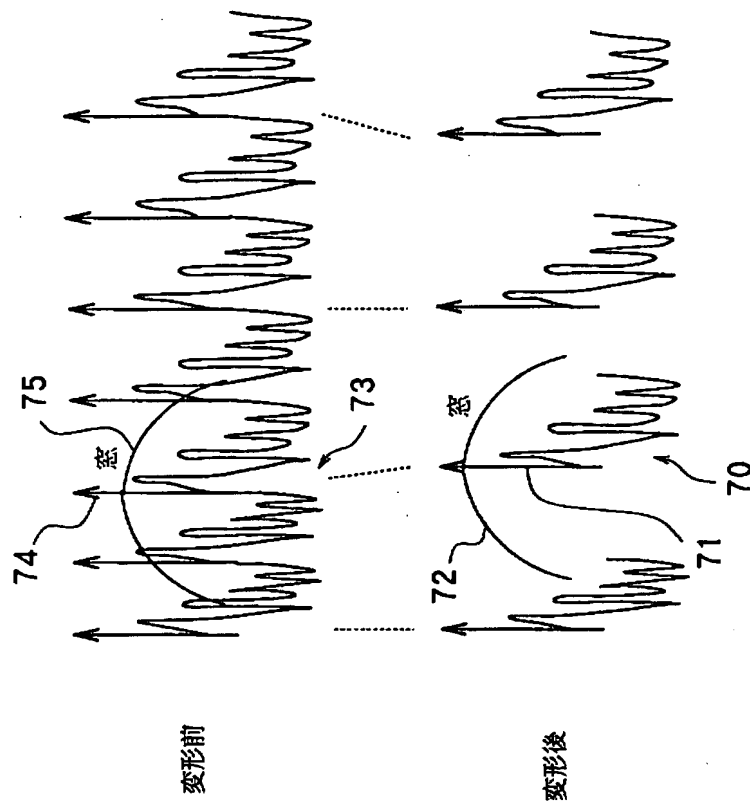
【図 5】



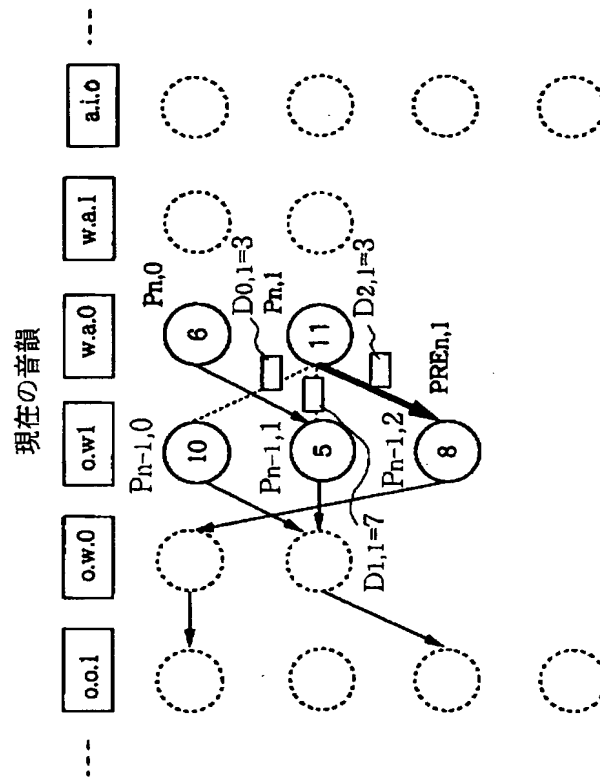
【図 6】



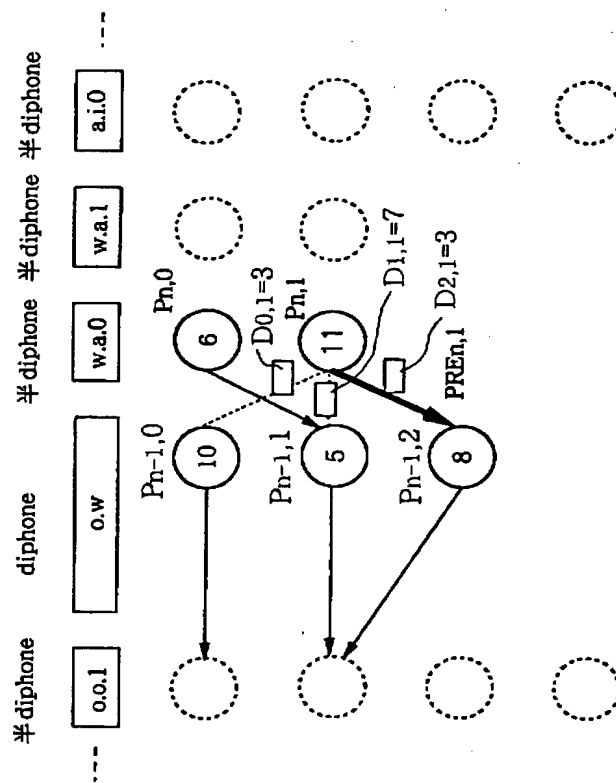
【図7】



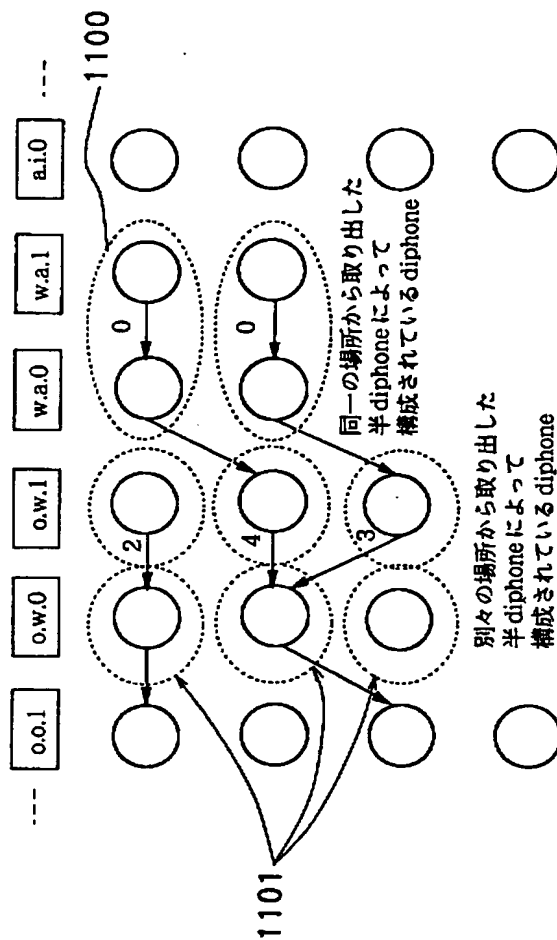
【図 8】



【図 9】



【図 10】



【図 1 1】

[illegible]

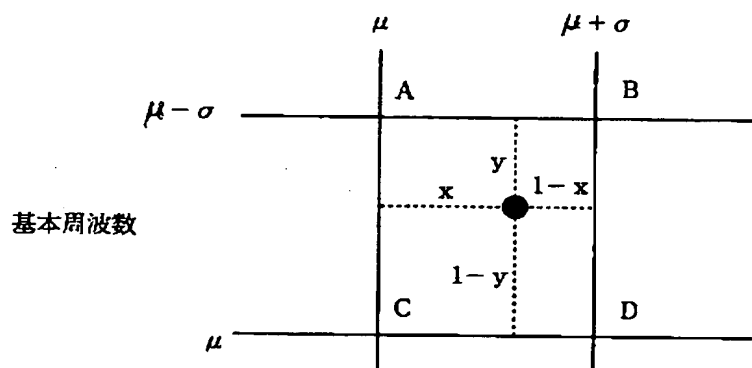
【図 1 2】

音韻時間長

	$\mu - 2\sigma$	$\mu - \sigma$	μ	$\mu + \sigma$	$\mu + 2\sigma$
$\mu - 2\sigma$	26	25	15	13	14
$\mu - \sigma$	20	24	19	18	17
μ	15	19	10	13	20
$\mu + \sigma$	19	22	16	22	26
$\mu + 2\sigma$	25	27	26	30	31

【図 1 3】

音韻時間長



【書類名】 要約書

【要約】

【課題】 接続歪と変形歪に基づく歪が小さくなるように音声素片を選択して、音声合成の音質劣化を抑制する。

【解決手段】 所定の音韻環境に対応付けて複数の音声素片を保持する音声素片保持部 2 0 6 から音韻環境に対応する複数の音声素片を抽出し (S 4 0 2)、それら抽出された複数の音声素片のそれぞれの歪を算出し (S 4 0 4)、音韻環境に基づいて決定される所定区間内で最小歪を求め (S 4 0 5)、最小歪経路となる音声素片列を選択し (S 4 0 6)、その音声素片を編集・接続して音声合成を行なう。

【選択図】 図 4

出 願 人 履 歴 情 報

識別番号 [000001007]

1. 変更年月日 1990年 8月30日
[変更理由] 新規登録
住 所 東京都大田区下丸子3丁目30番2号
氏 名 キヤノン株式会社